



GLOBAL JOURNAL OF COMPUTER SCIENCE AND TECHNOLOGY: C  
SOFTWARE & DATA ENGINEERING

Volume 14 Issue 2 Version 1.0 Year 2014

Type: Double Blind Peer Reviewed International Research Journal

Publisher: Global Journals Inc. (USA)

Online ISSN: 0975-4172 & Print ISSN: 0975-4350

## Research of Data Mining Algorithm based on Cloud Database

By Tianxiang Zhu, Xia Zhang, Dan Zhang & Xin Liu

*Shenyang University of Technology, China*

*Abstract-* There is an immense amount of data in the cloud database and among these data, much potential and valuable knowledge are implicit. The key point is to discover and pick out the useful knowledge, and to do so automatically. In this paper, the data model of the cloud database is analyzed. Through analyzing and classifying, the common features of the data are extracted to form a feature data set. The relationships among different areas in the data are then analyzed, from which the new knowledge can be found. In the paper, the basic data mining model based on the cloud database is defined, and the discovery algorithm is presented.

*Keywords:* cloud database, data mining, association rules, classification characteristic.

*GJCST-C Classification:* H.2.8



*Strictly as per the compliance and regulations of:*



© 2014. Tianxiang Zhu, Xia Zhang, Dan Zhang & Xin Liu. This is a research/review paper, distributed under the terms of the Creative Commons Attribution-Noncommercial 3.0 Unported License <http://creativecommons.org/licenses/by-nc/3.0/>, permitting all non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

# Research of Data Mining Algorithm based on Cloud Database

Tianxiang Zhu <sup>α</sup>, Xia Zhang <sup>ο</sup>, Dan Zhang <sup>ρ</sup> & Xin Liu <sup>ω</sup>

**Abstract-** There is an immense amount of data in the cloud database and among these data, much potential and valuable knowledge are implicit. The key point is to discover and pick out the useful knowledge, and to do so automatically. In this paper, the data model of the cloud database is analyzed. Through analyzing and classifying, the common features of the data are extracted to form a feature data set. The relationships among different areas in the data are then analyzed, from which the new knowledge can be found. In the paper, the basic data mining model based on the cloud database is defined, and the discovery algorithm is presented.

**Keywords:** cloud database, data mining, association rules, classification characteristic.

## I. INTRODUCTION

Cloud computing is derived from technologies such as distributed processing, parallel processing, grid computing, etc. It is an emerging approach to sharing the infrastructure architecture[1]. It distributes all the computing tasks on the resource pool that is made of many computers, making sure all the application systems can acquire desired computing power, memory space and software service according to their demand[2]. All the computing is provided to the terminal user by the form of service, and all the application software in the cloud as shared resources. A cloud database is a database deployed and virtualized in the cloud computing environment. It is predicated that as it develops overtime, more and more people and companies will store all their data in the cloud, which will make data mining based on the cloud computing one of the trends in the future data mining systems[3].

There is a massive amount of data in the cloud database, and among them, lives potentially valuable knowledge. How to discover such useful knowledge is the key point in database research. Data mining is the process of picking out the hidden knowledge and regulations, which possess potential value that could influence decision making[4]. Data mining namely refers to the knowledge discovery from a database and is comprised of the following procedures: data pre-processing, data alternating, data mining operation, rule expression and evaluation[5]. A data mining system includes: control unit - used to control all parts in a harmonious way; database interface – used to generate

Author <sup>α</sup> <sup>σ</sup> <sup>ρ</sup>: School of Software, Shenyang University of Technology Shenyang, Liaoning, China. e-mail: zhutianxiang@gmail.com

Author <sup>ω</sup>: Liaoning Northern Laboratory CO., LTD, Shenyang, Liaoning, China.

and process data according to the given query; database - used to store and manage relevant knowledge; focus - refers to the data extent that needs to be inquired; model extracting - refers to the various data mining algorithms; and finally, knowledge evaluation- used to evaluate the extracted conclusion[6].

## II. CLOUD DATABASE

A distributed database is a logical set of the databases at various sites or nodes in a computer network and logically, such databases belong to the same system[7]. Different from the traditional distributed database, a cloud database contains isolated as well as shared data; a cloud database can be designed by using different data models, which mainly include the key-value model and relationship model.

All data of the key-value model, including the rows and columns, are stored in the cells of a table. Contents are partitioned by row, the rows make up a tablet, and the tablet is stored on a server node.

### a) Row Key

Data is maintained in the lexicographic order on the row key. For a table, a row interval is dynamically partitioned according to the value of the row key and is the basic unit in which load balancing and data distribution are performed. Row keys are distributed amongst data servers.

### b) Column Key

Column keys are grouped into sets of many “column families” and are the basic units in which access control is performed. All data stored in a column family usually belong to the same data type, which means data is compressed at a higher rate. Data can be stored in a column key of the column family.

### c) Timestamp

Each cell contains multiple versions of the same data and these versions are indexed by the timestamp. Data model for key-value cloud database is as shown in Fig. 1:

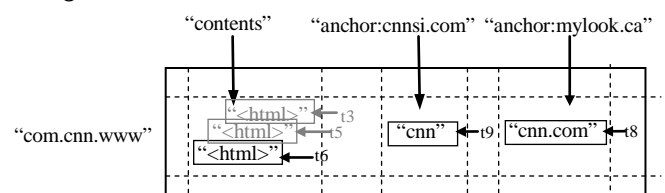


Figure 1 : Data model for key-value cloud database

The data model for the relational cloud database involves such relevant terms as row group and table group. A table is a logical relationship and includes a partitioning key, which is used for partitioning the table. The set of many tables with the same partitioning key is called a table group. In that table group, the set of rows with the same partitioning key value is called a row group. The rows in that row group are always allocated to the same data node. Each table group contains many row groups, which are allocated to different data nodes. A data partition contains many row groups, so each data node stores all rows with a certain partitioning key value. The data model for the relational cloud database is as shown in Fig. 2:

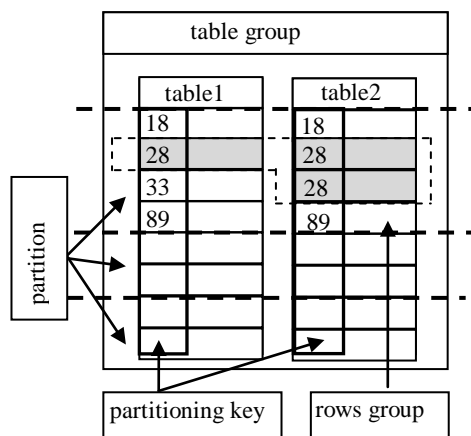


Figure 2 : Data model for relational cloud database

### III. DATA MINING FOR ASSOCIATION RULES

#### a) Model of the Association Rules

The normal target of the association rules is to discover the data relations among the data item set in the relationship type cloud database. Through mining based on the association rules, we can discover the relevance of the data.

In the subject item set, there are some target features in the relationship type cloud database. For instance, the commodity data item set in the commercial behavior analysis {T-shirt, coat, shoes, milk, bread ... }; data item set in the medical diagnosis analysis {hypertension, diabetes... }.

Classifying item set has the similar features with the subject item set, for instance, customer data item set in the commercial behavior analysis {vocation, gender, age... }; diagnosis behavior in medical diagnosis and signs and symptoms item set {smoking, polysaccharide, hyperlipidemia ... }.

Sample item set, which has both the features in the subject item set and the transaction data item set in the classifying item set. For instance, transaction data in the commercial activity analysis { {Zhangsan, milk}, {Zhangsan, bread}, {Lisi, T-shirt} ... }, health check

information in the medical diagnosis {{ Zhangsan, smoking, hypertension}, {Lisi, hyperlipidemia, diabetes}... }.

Through the mining based on the association rules, we can find that 90% of the customers who bought milk also bought bread; 50% of the patients who have hyperlipidemia also have diabetes.

The common targets of the association rules are transaction databases with the characters of subjects oriented item set. In practice, most databases are relational, and many applications and the required knowledge are from many different item sets (or multi-item set for simplicity) . For relational databases, it is difficult to describe the complicated association rules between the multi-item set with models of general association rules. We present the association rules model of multi-item set for the relational databases:

*Definition 1:* I is the subject item set, J is the taxonomy item set, each transaction corresponds to a subset T of the subject item sets and a taxonomy item U of the taxonomy item sets, called T belonging to class U.

*Model 1:* it is supposed that  $R=(r_1,r_2,\dots,r_n)$  is the rows group in the relational cloud database,  $r_k$  is one of the rows item set, D is a sample item set relevant to R, and each sample d corresponds to one rows item set, i.e.  $d \subseteq R$ . Each sample is marked with SID (sample identifier). As for the classifying item set X, only when  $X \subseteq d$ , the sample X belongs to d. association rules is a formula like  $X \subseteq d \Rightarrow Y \subseteq d$ , it can be  $X \Rightarrow Y$ , therein,  $X \subseteq R$ ,  $Y \subseteq R$  and  $X \cap Y = \Phi$ .

The rule  $X \Rightarrow Y$  in the sample item set D is constrained by degree of confidence C and degree of support S. Degree of confidence C is defined as C% in the transaction X in D also contains Y. Degree of support S is defined as transaction  $X \cup Y$  accounts for S% in D. Degree of confidence represents the strength of the rule, while Degree of support means the frequency of the model, which is shown in the rule.

In the cloud database containing cases information, 66% of the crime site in the theft cases happened in factories, so the C is 66%. Theft cases and factory cases account for 17% of the total cases, so the S is 17%.

The data frequency item set can be defined as the data item set where the degree of support S is over the pre-defined minimum degree of support S. The association rules with high degree of support S and degree of confidence C is considered strong association rules, otherwise it is considered weak association rules. Association rules mining means to find the line group that accord to the strong association rules in the database.

The procedure for mining these kinds of association rules of multi-item set is as follows:

1. Divide transaction D into several transaction subset  $D'=\{D1',D2',\dots,Dn'\}$  according to taxonomy item sets.
2. For all  $D1'<D'$  Do  
Find the strong sets of the main subject item  
Derive the association rules using the strong set
3. Next

These association rules of the multi-item set possess a feature where only one value is available in each sample (SID) set. With this method, mining the data's association rules is applicable for one-to-many relational databases. This is more practical and expands the mining range for the association rules.

In practice, most of the applications and knowledge is from the multiple data item set. For example, we regard a criminal case as the sample item set. For each case, there is one mark SID, several suspects, as well as several methods by which the crime is committed. So we can first take the education level of the suspects as one data item set, and the methods of committing crime as another data item set.

There are association rules with several multi-item sets, the association rules model can be termed as:

*Model 2:* It is supposed that  $I=(i_1,i_2,\dots,i_n)$  is a classifying item set,  $J=(j_1,j_2,\dots,j_m)$  is another one, D is a sample item set, each sample has two classifying item sets  $T(T\subseteq I)$  and  $T'(T'\subseteq J)$ , and each sample is marked with SID. The formula is  $X\subseteq I \Rightarrow Y\subseteq J$ , degree of confidence C can be termed as that in sample where D contains  $X\subseteq I$ , C% has  $Y\subseteq J$ , degree of support S can be defined as transaction with  $X\subseteq I$  and  $Y\subseteq J$  accounts for S% in D.

#### b) Mining Algorithm

There are many algorithms in the association rules, and the representative Apriori Algorithm follows the rule that the sub-item sets of all the strong item sets are classified to the strong item sets, while the super item sets of the weak item sets are weak item sets.

The first pass of the algorithm simply counts item occurrences to determine the strong 1-itemsets. A subsequent pass, pass k, consists of two phases. First, the strong item sets L found in the (k-1)th pass are used to generate the candidate item sets  $C_k$ , using the apriori-gen function. Next, the database is scanned and the support of candidates in  $C_k$  is counted. For fast counting, we need to efficiently determine the candidates in  $C_k$  that are contained in a given samples.

As for the association rules of multiple data item sets, we need to have strong item sets L1 with item 1, and then we can have C2 from L1 with the item 2, after this we can have L2, based on this method we can finally have  $C_k$ , and get Lk from the database.

Classifying item set D into m classifying item sets D1, D2, ... Dm according to the separating item set J, then we can find out the association rules after using Apriori Algorithm to each sub-sample item set D.

```

for(j=1;j<=m;j++) do
begin
Lj,1={large 1-items};
for(k=2;Lj,k-1≠Φ;k++) do
begin
Ck=apriori-gen(Lj,k-1);
forall samples s∈Dj do
begin
Cs=subset(Ck,s);
forall candidates c∈Cs do
c.count++;
end
Lj,k={c∈Ck|c.count>=minsup}
end
Answer=Uj,kLj,k;
end;

```

$L_{j,1}$  represents the strong item set in  $D_j$  sub sample item set, which will generate K item in  $D_j$ , scan the database to have  $L_{j,k}$ , we finally can have  $D_1, D_2, \dots, D_m$  strong item set from the sub sample item set.

Since Model 2 corresponds to two classifying item sets and each sample  $S\subseteq D$  includes classifying item set I and J, 1-itemsets represent the strong item sets we select from I and J, which is  $L_{i,j}$ . From  $L_{i,j}$  we can have  $C_{1,2}$  from  $L_{1,2}$ , done with the similar manner, and then get  $L_{1,k}$ . From  $L_{1,1}$ , we can have  $C_{2,1}$  from  $L_{2,1}$ , the algorithm is:

```

L1,1={1-itemsets x∈I, y∈J};
for(i=2;i=n;i++) do
begin
for(j=2;j<=m;j++) do
begin
Ci,j=apriori-gen(Li,j-1);
forall samples s∈D do
begin
Cs=subset(Ci,j,s);
forall candidates c∈Cs do
c.count++;
end
Li,j={c∈Ci,j|c.count>=minsup}
end
Answer=Ui,jLi,j;
End

```

In management information systems, the relational database is widely used; the connection among different data is one-to-many and many-to-many, so it is universal to discover knowledge in the database. As the cloud age is coming, data mining from the cloud data is more practical. The mining method that is used in the association rules is applied to the cloud database, making the association rules more

practical and universal. This paper also extends the Apriori Algorithm into association rules mining model, which realize the mining multi-item set association rules.

#### IV. DATA MINING FOR CLASSIFICATION CHARACTERISTIC RULE

Knowledge discovered from a database with massive data is diversified in variety. Knowledge classification refers to clustering or classifying tuples in the database to divide these tuples into different categories by characteristic rules extracted from a certain target class, and thus achieve the purpose of describing the characteristics of the tuples of that class.

Clustering refers to categorizing a group of individuals into several categories, which means those with the same characteristic are classified as one category. Clustering is a process in which a data object with multiple attributes is continuously classified. In such process, classification is automatically executed by the classification algorithm to divide the data into several classes by identifying data features. A relational database mainly containing character information may be equivalently partitioned into equivalence classes according to the concept of equivalence class. The resulting equivalence classes are a group of classes. The characteristics of each class are further analyzed and this can lead to the determination of the classification characteristic rules. Such analysis process is of practical significance. For example, symptoms and reaction characteristics of various diseases can be determined by analyzing a great amount of medical diagnosis cases.

##### a) Classification Model For Key-Value Model Based Cloud Database

Let D be a key-value model based cloud database, K represents the set of all row keys in D with the formula  $K = \{k_1, k_2, \dots, k_n\}$ , At represents the set of all column keys in D with the formula  $A = \{a_1, a_2, \dots, a_m\}$ , V represents the dataset of certain attribute characters of the column keys with the formula  $V = \{v_{11}, v_{12}, \dots, v_{mn}\}$  and f represents a function of a and k with the formula  $V_{i,j} = f(a_i, k_j)$ .

**Definition 2:** For  $\forall a \in A_i$  ( $A_i$  is the dataset of column keys,  $A_i \subseteq A$ ), if  $k_i \in K, k_j \in K, i \neq j$  and  $f(a, k_i) = f(a, k_j)$ , then  $k_i$  is said to be equivalent to  $k_j$  based on the dataset of column key attributes At, and the set of all equivalent row keys based A is called equivalence class based on the dataset of column key attributes A; all row sets in K are classified by equivalence class and the classification result is called A-based classification:  $K = \{K_1, K_2, \dots\}$ .

**Definition 3:**  $K = \{K_1, K_2, \dots\}$ , K is the At-based classification and a column key in At is called a classification. The attribute value of the column key in At is called the name of the classification.

**Definition 4:** Let D be a key-value model based cloud database,  $S_k$  represents the amount of the latest timestamps in the row key set, At is a column key set of D, Y is a At-based equivalence class and  $S_y$  is the amount of the latest timestamps in the set of row keys in Y, then  $S = S_y / S_k * 100\%$  is said to be the classification support degree of the equivalence class Y.

**Example 1:** Let D be a key-value model based cloud database, K is the set of all row keys in D, A is the set of all column keys in D and At is a subset of A. V is the dataset of certain attribute characters of the column keys and each data has the latest timestamp.

$$K = \{k_1, k_2, k_3, k_4, k_5, k_6, k_7, k_8, k_9, k_{10}, k_{11}, k_{12}\}$$

$$A = \{a_1, a_2, a_3\}$$

$$At = \{a_1\}$$

$$V = \{v_{10}, v_{11}, v_{12}, v_{20}, v_{21}, v_{30}, v_{31}, v_{32}\}$$

The values of  $f(k,a)$  are as shown in Table 1.

Table 1 : The values of  $f(k,a)$

K \ A	A1	A2	A3
k <sub>1</sub>	v <sub>11</sub>	v <sub>20</sub>	v <sub>32</sub>
k <sub>2</sub>	v <sub>10</sub>	v <sub>21</sub>	v <sub>30</sub>
k <sub>3</sub>	v <sub>12</sub>	v <sub>20</sub>	v <sub>30</sub>
k <sub>4</sub>	v <sub>11</sub>	v <sub>21</sub>	v <sub>30</sub>
k <sub>5</sub>	v <sub>11</sub>	v <sub>20</sub>	v <sub>32</sub>
k <sub>6</sub>	v <sub>12</sub>	v <sub>20</sub>	v <sub>30</sub>
k <sub>7</sub>	v <sub>10</sub>	v <sub>21</sub>	v <sub>31</sub>
k <sub>8</sub>	v <sub>11</sub>	v <sub>21</sub>	v <sub>31</sub>
k <sub>9</sub>	v <sub>11</sub>	v <sub>20</sub>	v <sub>32</sub>
k <sub>10</sub>	v <sub>10</sub>	v <sub>21</sub>	v <sub>31</sub>
k <sub>11</sub>	v <sub>11</sub>	v <sub>20</sub>	v <sub>32</sub>
k <sub>12</sub>	v <sub>10</sub>	v <sub>21</sub>	v <sub>31</sub>

In the above mentioned database, the field {a1} in the column key set At can be classified into three classes:

$K_1 = \{k_1, k_4, k_5, k_8, k_9, k_{11}\}$ ;  $K_2 = \{k_2, k_7, k_{10}, k_{12}\}$ ;  $K_3 = \{k_3, k_6\}$ .  $\{K_1, K_2, K_3\}$  is a class based on the column key set At, the name of that class is {v11, v10, v12} and its classification support degree is {50%, 33.33%, 16.67%}.

##### b) Classification Model for Relational Model Based Cloud Database

Let D be a relational model based cloud database and T be a table group of D, P represents the set of partitioning keys in T with the formula  $P = \{p_1, p_2, \dots, p_n\}$  and R represents the set of row groups of the partitioning key Pi with the formula  $R = \{r_1, r_2, \dots, r_n\}$ .

**Definition 5:** R represents the set of row groups of the partitioning key Pi with the formula of  $R = \{r_1, r_2, \dots, r_n\}$

and is called a class based on the partitioning key  $P = \{p_1, p_2, \dots, p_n\}$ .  $p_n$  is called the name of the class  $r_n$ .

*Definition 6:* Let  $D$  be a relational model based cloud database and  $T$  be a table group of  $D$ ,  $S_i$  represents the record count of all row groups, the row group set  $R$  is the class based on the partitioning key  $P$  and  $S_y$  represents the record count of  $Y$  row groups in  $R$ , then  $S = S_y/S_i * 100\%$  is said to be the classification support degree of the class  $Y$ .

*Example 2:* Let  $D$  be a relational model based cloud database and  $T$  be a table group of  $D$ ,  $P$  is the partitioning key and the value of  $P$  is  $\{p_1, p_2, p_3\}$ , then the corresponding row group is  $\{r_1, r_2, r_3\}$ , namely:

$P = \{p_1, p_2, p_3\}$   
 $R = \{r_1, r_2, r_3\}$   
 $r_1 = \{v_{11}, v_{12}, v_{13}, v_{14}, v_{15}, v_{16}\}$   
 $r_2 = \{v_{21}, v_{22}, v_{23}\}$   
 $r_3 = \{v_{31}\}$

The above mentioned database can be partitioned into three classes based on the partitioning key  $P$ :

$r_1 = \{v_{11}, v_{12}, v_{13}, v_{14}, v_{15}, v_{16}\}$ ,  $r_2 = \{v_{21}, v_{22}, v_{23}\}$ ,  $r_3 = \{v_{31}\}$ , the name of the class is  $\{r_1, r_2, r_3\}$  and the classification support degree is  $\{60\%, 30\%, 10\%\}$ .

For the cloud database  $D$ , all classification support degrees  $\{S_1, S_2, S_3, \dots\}$  can be obtained according to a certain classification  $R = \{R_1, R_2, R_3, \dots\}$ .

*Definition 7:* Let  $S_p$  be a given threshold,  $0 \leq S_p \leq 1$ . Those classes with a classification support degree  $S \geq S_p$  are called strong class and those with a classification support degree  $S < S_p$  are called weak classes.

In mining knowledge from massive data, we usually are concerned about and interested in data classes with higher classification support degree, namely the strong classes. Strong classes contain more representative knowledge.

*c) Classification Characteristic Rule Model*

According to the definitions as mentioned above, data in a database can be classified and the characteristics of the strong classes need to be further analyzed.

*Definition 8:* Let  $E$  be a  $A_t$ -based class,  $A_i$  is the complementary set of  $A_t$  against the attribute  $A$ ,  $A_i \subseteq A$ ,  $B$  is the subset of  $A_i$ , the equivalence class  $T$  based on  $B$  is called the characteristic domain in the class  $E$ , and the value  $\{b_1, b_2, \dots\}$  of the attribute in  $B$  is called the characteristics in the class  $E$ .

*Definition 9:* Let  $e_c$  be the record count of the class  $E$ , and  $t_c$  be the record count of the characteristic domain  $T$ , then  $C = t_c/e_c * 100\%$  is said to be the confidence degree of characteristic.

*Definition 10:* Let  $C_p$  be a given threshold,  $0 \leq C_p \leq 1$ , a characteristic domain with the confidence degree of

characteristic  $C \geq C_p$  is called a strong characteristic domain and a characteristic domain with the confidence degree of characteristic  $C < C_p$  is called a weak characteristic domain. The value of the field with strong characteristic domain is called a strong characteristic, while the value of the field with weak characteristic domain is called a weak characteristic.

Strong characteristics in a strong class are usually representative knowledge and can be expressed as:

$(E, T, C_p)$   
 $E$ : class  
 $T$ : characteristic  
 $C_p$ : confidence degree of characteristic

Discovery Algorithm for Classification Characteristic Rules.

$D$  is a cloud database and  $A$  is the set of classification attributes of  $D$ .

```

For all  $A_i \subseteq A$  Do
    Obtain the set of  $A_t$ -based classes  $E_1, \dots, E_m$ ;
    Obtain the classification support degree of the class set  $E_1, \dots, E_m$   $S = \{S_1, \dots, S_m\}$ ;
    For  $i=1$  To  $m$  Do
        If  $S_i \geq S_p$  Then
            Obtain all characteristic domains  $T_1, T_2, \dots, T_k$ ;
            Obtain the confidence degree of characteristic of the characteristic domain set  $T = \{T_1, T_2, \dots, T_k\}$ 
             $C = \{C_1, C_2, \dots, C_k\}$ ;
            For  $j=1$  To  $k$  Do
                if  $C_j \geq C_p$  Then
                     $(E_i, T_j, C_j) \Rightarrow$  result base
            Endif
        Next
    Endif
Next
Next
Next
    
```

V. APPLICATION OF THE CLASSIFICATION CHARACTERISTIC RULES IN CASE INFORMATION SYSTEMS

Suppose the related property is the case type, selected site and the way of commit, the related degree of the door smashed versus picked is 0.8, the given threshold of the related degree is 2.5, two cases as shown in Table 2:

Table 2: These Two Cases are Related

	Case kind	Selected site	Way of commit
Case1	Burglary	residence	door picked
Case2	Burglary	residence	door smashed

Based on the above definitions, as long as the related degrees of the related properties are known, the

related cases can be discovered. The values of the related degrees are provided by the field experts according to the field knowledge. In order to express the related degrees, a related degree matrix  $M_a$  is defined as follows:

$$M_a = \begin{bmatrix} C_{11} & C_{12} & \dots & C_{1n} \\ C_{21} & C_{22} & \dots & C_{2n} \\ \dots & \dots & \dots & \dots \\ C_{m1} & C_{m2} & \dots & C_{mn} \end{bmatrix} \quad (1)$$

$C_{ij}$ : related degree of element  $j$  to element  $i$  of property  $a$ .

$M_a$  is a symmetrical matrix, so only consider the lower triangle.

The related degree matrix of the way of commit is as shown in Table 3:

Table 3: The Related Degree Matrix of the Way of Commit

	Pick door lock	Tag along after	Smash door	Cut bag	Pretend to be conceal	Conceal
Prize door lock	1	0	0	0	0	0
Tag along after	0	1	0	0	0	0
Smash door	0.8	0.5	1	0	0	0
Cut bag	0	0.9	0	1	0	0
Pretend to be	0	0.8	0.3	0	1	0
Conceal	0	0	0	0	0	1

During case analysis, the related degree matrix of all related properties must be known. The sum of case related degrees can be found based on the related properties and the related degree matrix and all the related cases can be found based on the given threshold  $C_p$ .

Input  $U$ , an information system, has  $n$  records,  $M_a$  is a symmetrical matrix,  $M_{ak}[i,j]$  seeks the correlation degree in the correlation matrix.

Output  $L$ , a set of case related to it.

For  $i=1$  To  $n-1$  ( $n$  is the record number)

For  $j=i+1$  TO  $n$

$A[i, j] = 0$

For  $k=1$  TO  $m$

$A[i,j]=A[i,j]+M_{ak}[i,j]$

Next

If  $A[i,j] \geq C_p$

$uj \cup \{u \mid u \text{ relates to } u_i\}$

End If

Nexts

$L = \cup \{u \mid u \text{ relates to } u_i\}$

Next

Output  $L$

Following the idea of converging classes, the case information is divided into many kinds with the equivalent dividing method. Define the base value of kinds as classifying the support degrees. The kinds can be divided into strong class ones and weak class ones. The weak class has too small classifying support degrees, no practical meanings and can be neglected. For the strong class, Rough set theory can be used to analyze their common features and form the classifying characteristic regulations.

Data was mined from the database for the experimental criminal case information system by using the aforementioned algorithm. Taking the crime approach table group as an example, the table contains 100762 records. Given a classification support degree threshold of 10%, and a characteristic confidence degree threshold of 20%, 359 classification characteristic rules were mined, for example:

(Residential house, night, 23.4%)

(Rubbery, less than RMB10000, 93.3%)

## VI. CONCLUSIONS

The data mining technique is new to the information society. Many subjects need to be studied in this field. In many professions, a certain amount of databases have been accumulated, in which some hidden knowledge needs to be discovered. Starting with the concept of set theory, the data model for the cloud database was analyzed; the model and algorithm for mining classification characteristic rules from cloud database were designed to make data mining of classification characteristic rule more practical.

The abstracted related knowledge models presented in this paper can be put into practice in the public security affairs, such as case chaining, which is one of the highly demanded, complex tasks in the public security affairs. The presented methods about the related case data mining in this paper promote the work effect of the chained cases. On the case material analysis, the mining of the classifying characteristic regulations help users with their classifying work and overcome the weaknesses that exist in the old statistics method, in which repeated experimentation are required.

## REFERENCES RÉFÉRENCES REFERENCIAS

1. Lin ZY, Lai YX, Lin C, Xie Y, Zou Q. Research on cloud databases. Journal of Software, 2012,23(5):124-137.
2. Xu M, Gao D, Deng C, Luo ZG, Sun SL. Cloud computing boosts business intelligence of

- telecommunication industry. In: Jaatun MG, Zhao GS, Rong CM, eds. Proc. of the 1st Int'l Conf. on Cloud Computing (CloudCom 2009). Berlin: Springer-Verlag, 2009. 224–231.
3. Feng DG, Zhang M, Zhang Y, Xu Z. Study on cloud computing security. *Journal of Software*, 2011,22(1):71–83.
  4. Zhu TX, W P, Zhang D. Application of the data mining technique in case information systems. *ICCSEE 2012*, 3(1):43-46.
  5. Zhu TX, Li L, Xu ZW, Technology for Mining Classification-Characteristic Rules, *Journal of Shenyang Polytechnic University*, 1999.(x)x:22-24 .
  6. He JJ, Ye CM, Wang XB, Huang ZX, Liu QL. Cloud Computing-Oriented Data Mining System Architecture, *Application Researche of Computers*, 2011,28(4):1372-1374.
  7. Abouzeid A, Bajda-Pawlikowski K, Abadi DJ, Rasin A, Silberschatz A. HadoopDB: An architectural hybrid of MapReduce and DBMS technologies for analytical workloads. *PVLDB*, 2009,2(1):922–933.





This page is intentionally left blank